

Power Expectation Propagation for the Relevance Tagging Machine

Dmitry Molchanov
Dmitry Kondrashkin
Dmitry Vetrov

November 06, 2015

Complex distributions

$$p(\mathbf{X}, \mathbf{Y}, \mathbf{Z})$$

$$p(\mathbf{X}) = \int p(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) d\mathbf{Y} d\mathbf{Z} = ?$$

$$\mathbb{E}_{p(\mathbf{X})} f(\mathbf{X}) = \int f(\mathbf{X}) p(\mathbf{X}, \mathbf{Y}, \mathbf{Z}) d\mathbf{X} d\mathbf{Y} d\mathbf{Z} = ?$$

Distribution approximation

Consider a distribution $p(x)$ from a complex distribution family. The goal is to approximate it with a distribution $q(x)$ from a simple distribution family:

$$p(x) \approx q(x)$$

$$p(x) = \prod_{i=1}^n t_i(x), \quad t_i \in \mathcal{F}_i - \text{complex family of distributions}$$

$$q(x) = \prod_{i=1}^n \tilde{t}_i(x), \quad \tilde{t}_i \in \mathcal{Q}_i - \text{simple family of distributions}$$

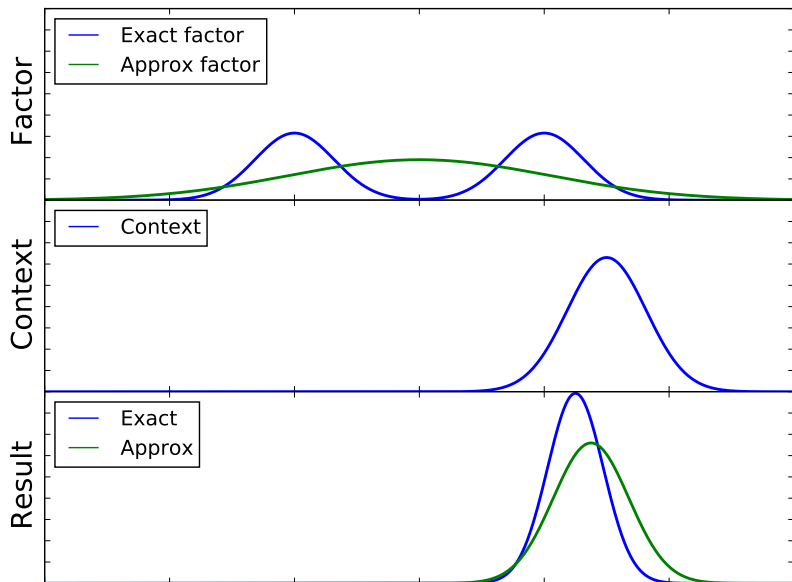
Distribution approximation: naive approach

$$p(x) = \prod_{i=1}^n t_i(x), \quad q(x) = \prod_{i=1}^n \tilde{t}_i(x)$$

$$\tilde{t}_i(x) \approx t_i(x)?$$

Bad idea!

Distribution approximation: naive approach



Contextual approximation

$$p(x) = \prod_{i=1}^n t_i(x), \quad q(x) = \prod_{i=1}^n \tilde{t}_i(x)$$

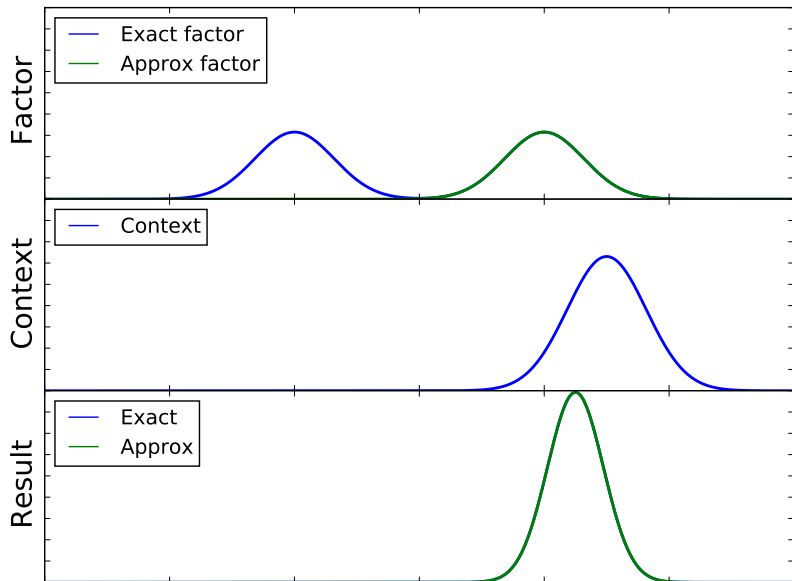
Context for i -th factor:

$$q^{(i)}(x) = \frac{q(x)}{\tilde{t}_i(x)}$$

Contextual approximation:

$$q^{(i)}(x)\tilde{t}_i(x) \approx q^{(i)}(x)t_i(x)$$

Contextual approximation



KL-divergence

A brief reminder: Kullback–Leibler divergence

$$\text{KL}(p\|q) = \int p(x) \log \frac{p(x)}{q(x)} dx$$

Some properties:

- $\text{KL}(p\|q) \geq 0$
- $\text{KL}(p\|q) = 0 \Leftrightarrow p(x) = q(x)$ almost everywhere
- Can be used to measure the difference between p and q , ...
- ... but is not a metric and is not symmetric

KL-projection

Contextual approximation:

$$q^{i}(x)\tilde{t}_i(x) \approx q^{i}(x)t_i(x)$$

KL-projection:

$$q^{i}(x)\tilde{t}_i(x) = \arg \min_{g \in \mathcal{Q}} \text{KL}(q^{i} \cdot t_i \parallel g) = \text{proj}_{\mathcal{Q}}[q^{i} \cdot t_i]$$

KL-projection in exponential family

Let \mathcal{Q} be a subset of exponential family with parameters $\boldsymbol{\theta}$:

$$q(\mathbf{x} \mid \boldsymbol{\theta}) = \frac{g(\mathbf{x})}{h(\boldsymbol{\theta})} \exp(\boldsymbol{\theta}^T \boldsymbol{\phi}(\mathbf{x}))$$

$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_m(\mathbf{x}))^T$ are sufficient statistics of \mathcal{Q} .

In this case KL-projection on \mathcal{Q} is equivalent to moment matching:

$$g = \text{proj}_{\mathcal{Q}}[f]$$

$$\Leftrightarrow$$

$$\int \phi_j(\mathbf{x}) g(\mathbf{x}) d\mathbf{x} = \int \phi_j(\mathbf{x}) f(\mathbf{x}) d\mathbf{x}, \quad \forall j = 1..m$$

EP algorithm

- 1 Initialize $q = \prod_{i=1}^n \tilde{t}_i$
 - 2 Recalculate context for all factors: $q^{i}(x) = q(x)/\tilde{t}_i(x)$
 - 3 Update $\tilde{t}_i = \text{proj}_{\mathcal{Q}}[q^{i} \cdot t_i]/q^{i}$
 - 4 Repeat steps 2-3 until convergence
- + Good in practice: fast and accurate
- Poor theoretical evidence to its convergence or approximation accuracy
- Fails if moment-matching in proj operator is intractable

Power EP motivation

EP fails if $\int \phi_j(x) q^{i_j}(x) t_i(x) dx$ is intractable.

However, $\int \phi_j(x) q^{i_j}(x) (t_i(x))^{\eta_i} dx$ may be tractable for some η_i .

For example:

$q(x)$ is a Gaussian, $\phi_1(x) = x$, $\phi_2(x) = x^2$

$t_i(x)$ is Student's t -distribution, $t_i(x) = \frac{1}{x^2 + 1}$

$$\int x^k \frac{1}{x^2 + 1} \mathcal{N}(x | \mu, \sigma^2) dx = ?$$

$$\int x^k \left(\frac{1}{x^2 + 1} \right)^{-1} \mathcal{N}(x | \mu, \sigma^2) dx = \int x^k (x^2 + 1) \mathcal{N}(x | \mu, \sigma^2) dx$$

Power EP

$$\text{EP: } \frac{q(x)}{\tilde{t}_i(x)} \tilde{t}_i(x) \approx \frac{q(x)}{\tilde{t}_i(x)} t_i(x)$$

$$\text{Power EP: } \frac{q(x)}{(\tilde{t}_i(x))^{\eta_i}} (\tilde{t}_i(x))^{\eta_i} \approx \frac{q(x)}{(\tilde{t}_i(x))^{\eta_i}} (t_i(x))^{\eta_i}$$

$$q^{\setminus i}(x) = \frac{q(x)}{(\tilde{t}_i(x))^{\eta_i}}$$

$$\tilde{t}_i = \left(\frac{\text{proj}_{\mathcal{Q}}[q^{\setminus i} \cdot t_i^{\eta_i}]}{q^{\setminus i}} \right)^{\frac{1}{\eta_i}}$$

+ Is applicable to a wider variety of models

Relevance Tagging Machine

We address the binary classification problem with binary features.

- $(\mathbf{x}_i, c_i)_{i=1}^n$ is the training set
- $\mathbf{x}_i \in \{0, 1\}^d$ is an object and $c_i \in \{0, 1\}$ is its class label
- $x_{ij} = 1 \Leftrightarrow \mathbf{x}_i$ is tagged by the tag j
- All tags affect the class label independently

Probabilistic model of the RTM:

$$\theta_j = P(c = 1 \mid x_j = 1), \quad P(c = 1 \mid \mathbf{x}, \boldsymbol{\theta}) = \frac{\prod_{j=1}^d \theta_j^{x_j}}{\prod_{j=1}^d \theta_j^{x_j} + \prod_{j=1}^d (1 - \theta_j)^{x_j}}$$

ARD

Bayesian *automatic relevance determination* (ARD) approach.

- Parameters are given independent priors
- Hyperparameters are trained by maximizing the evidence

Symmetrical Beta distribution:

$$\theta_j \sim \text{Beta}(\theta_j | \alpha_j, \alpha_j), \quad \alpha_j \in [1, +\infty)$$

- $\alpha_j = 1 \Rightarrow \theta_j^{MAP} = \theta_j^{ML}$
- $\alpha_j = +\infty \Rightarrow \theta_j = 0.5 \Rightarrow \theta_j$ is removed from the model

Evidence maximization:

$$\int P(\mathbf{c} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\alpha}) d\boldsymbol{\theta} \rightarrow \max_{\boldsymbol{\alpha}}$$

Likelihood approximation

$$P(\mathbf{c} | \mathbf{X}, \boldsymbol{\theta}) = \prod_{i=1}^n P(c_i | \mathbf{x}_i, \boldsymbol{\theta}) = \prod_{i=1}^n t_i(\boldsymbol{\theta}) \approx \frac{1}{Z} q(\boldsymbol{\theta})$$

$$q(\boldsymbol{\theta}) = \prod_{i=1}^n \tilde{t}_i(\boldsymbol{\theta}) = \prod_{i=1}^n \prod_{j=1}^d \text{Beta}(\theta_j | a_{ij}, b_{ij})$$

$$t_i(\boldsymbol{\theta}) = \frac{\prod_{j=1}^d \theta_j^{c_i x_{ij}} (1 - \theta_j)^{(1-c_i)x_{ij}}}{\prod_{j=1}^d \theta_j^{x_{ij}} + \prod_{j=1}^d (1 - \theta_j)^{x_{ij}}}$$

Likelihood approximation

$$\begin{aligned}t_i(\boldsymbol{\theta})^{-1} &= \left(\frac{\prod_{j=1}^d \theta_j^{c_i x_{ij}} (1 - \theta_j)^{(1-c_i)x_{ij}}}{\prod_{j=1}^d \theta_j^{x_{ij}} + \prod_{j=1}^d (1 - \theta_j)^{x_{ij}}} \right)^{-1} = \\&= \prod_{j=1}^d \theta_j^{x_{ij}(1-c_i)} (1 - \theta_j)^{x_{ij}(c_i-1)} + \\&\quad + \prod_{j=1}^d \theta_j^{-x_{ij}c_i} (1 - \theta_j)^{x_{ij}c_i}\end{aligned}$$

Power EP update for the RTM

$$\tilde{t}_i(\boldsymbol{\theta})^{new} = \left(\frac{\text{proj}[t_i(\boldsymbol{\theta})^{-1} \cdot q^i(\boldsymbol{\theta})]}{q^i(\boldsymbol{\theta})} \right)^{-1}$$

$$t_i(\boldsymbol{\theta})^{-1} \cdot q^i(\boldsymbol{\theta}) \propto \dots$$

$$\dots \propto \prod_{j=1}^d \theta_j^{A_{ij}^1 - 1} (1 - \theta_j)^{B_{ij}^1 - 1} + \prod_{j=1}^d \theta_j^{A_{ij}^2 - 1} (1 - \theta_j)^{B_{ij}^2 - 1}$$

Sufficient statistics of beta distribution $p(x) = \text{Beta}(x | a, b)$:

$$\phi(x) = (\log x, \log(1 - x))^T$$

Their expectation ($\psi(x)$ is the digamma function):

$$\mathbb{E}_{p(x)} \log x = \psi(a) - \psi(a + b),$$

$$\mathbb{E}_{p(x)} \log(1 - x) = \psi(b) - \psi(a + b)$$

Power EP update for the RTM

Moment matching:

$$\text{proj}[t_i(\boldsymbol{\theta})^{-1} \cdot q^{i}(\boldsymbol{\theta})] = \prod_{j=1}^d \text{Beta}(\theta_j \mid \tilde{A}_{ij}, \tilde{B}_{ij})$$

$$\begin{cases} \psi(\tilde{A}_{ij}) - \psi(\tilde{A}_{ij} + \tilde{B}_{ij}) = c_1 \\ \psi(\tilde{B}_{ij}) - \psi(\tilde{A}_{ij} + \tilde{B}_{ij}) = c_2 \end{cases}$$

It can be solved efficiently with 3–5 iterations of Newton's method [Minka, Estimating a Dirichlet distribution, 2012].

Convergence tricks

Power EP algorithm for the RTM:

$$q^{i+1}(x) = q(x) \cdot \tilde{t}_i(x)$$

$$\tilde{t}_i = \left(\frac{\text{proj}_{\mathcal{Q}}[q^{i+1} \cdot t_i^{-1}]}{q^{i+1}} \right)^{-1}$$

The algorithm, written in this form, does not converge.

We used several tricks to improve its stability:

- 1 Damping
- 2 Sequential update of factors
- 3 More damping
- 4 Addition of fixed prior distribution factors

Convergence tricks

Damping

- 1 **for** $k \leftarrow 1$ **to** *maximum iteration number or until convergence* **do**
- 2
$$\begin{aligned} \tilde{t}_i^{old} &:= \tilde{t}_i \quad \forall i \quad q^{i}(x) := q^{new}(x) \cdot \tilde{t}_i(x) \quad \forall i \\ \tilde{t}_i^{new} &:= \left(\frac{\text{proj}_{\mathcal{Q}}[q^{i} \cdot t_i^{-1}]}{q^{i}} \right)^{-1} \quad \forall i \\ \tilde{t}_i &:= \left(\tilde{t}_i^{old} \right)^{\gamma_i} \cdot \left(\tilde{t}_i^{new} \right)^{1-\gamma_i}, \quad \forall i, \gamma_i \in [0, 1) \end{aligned}$$
- 3 **end**

Parameters of \tilde{t}_i can be computed as a convex combination of old and new ones.

Convergence tricks

Sequential update of factors

```
1 for  $k \leftarrow 1$  to maximum iteration number or until  
   convergence do  
2    $\tilde{t}_i^{old} := \tilde{t}_i \quad \forall i$   
3   for  $i \leftarrow 1$  to  $n$  do  
4      $q^{new}(x) := \prod_{l=1}^{i-1} \tilde{t}_l^{new} \cdot \prod_{l=i}^n \tilde{t}_l$   
      $q^{\setminus i}(x) := q^{new}(x) \cdot \tilde{t}_i(x)$   
      $\tilde{t}_i^{new} := \left( \frac{\text{proj}_{\mathcal{Q}}[q^{\setminus i} \cdot t_i^{-1}]}{q^{\setminus i}} \right)^{-1}$   
5   end  
6    $\tilde{t}_i := \left( \tilde{t}_i^{old} \right)^{\gamma_i} \cdot \left( \tilde{t}_i^{new} \right)^{1-\gamma_i}, \quad \gamma_i \in [0, 1) \quad \forall i$   
7 end
```

Convergence tricks

More damping

1 **for** $k \leftarrow 1$ **to** *maximum iteration number or until convergence* **do**

2 $\tilde{t}_i^{old} := \tilde{t}_i \quad \forall i$

3 **for** $i \leftarrow 1$ **to** n **do**

4 $q^{new}(x) := \prod_{l=1}^{i-1} \tilde{t}_l^{new} \cdot \prod_{l=i}^n \tilde{t}_l$

$q^{\setminus i}(x) := q^{new}(x) \cdot \tilde{t}_i(x)$

$\tilde{t}_i^{new} := \tilde{t}_i^{\rho_i} \cdot \left(\left(\frac{\text{proj}_{\mathcal{Q}}[q^{\setminus i} \cdot t_i^{-1}]}{q^{\setminus i}} \right)^{-1} \right)^{(1-\rho_i)}$

5 **end**

6 $\tilde{t}_i := \left(\tilde{t}_i^{old} \right)^{\gamma_i} \cdot \left(\tilde{t}_i^{new} \right)^{1-\gamma_i}, \quad \gamma_i \in [0, 1) \quad \forall i$

7 **end**

Convergence tricks

Addition of fixed prior distribution factors

$$p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) = \prod_{j=1}^d \text{Beta}(\theta_j | \alpha_j^0, \alpha_j^0)$$

$$P(\mathbf{c} | \mathbf{X}, \boldsymbol{\theta}) p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) = p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) \prod_{i=1}^n P(c_i | \mathbf{x}_i, \boldsymbol{\theta}) =$$

$$= p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) \prod_{i=1}^n t_i(\boldsymbol{\theta}) \approx \frac{1}{Z} q(\boldsymbol{\theta})$$

$$q(\boldsymbol{\theta}) = p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) \prod_{i=1}^n \tilde{t}_i(\boldsymbol{\theta}) =$$

$$= p(\boldsymbol{\theta} | \boldsymbol{\alpha}^0) \prod_{i=1}^n \prod_{j=1}^d \text{Beta}(\theta_j | a_{ij}, b_{ij})$$

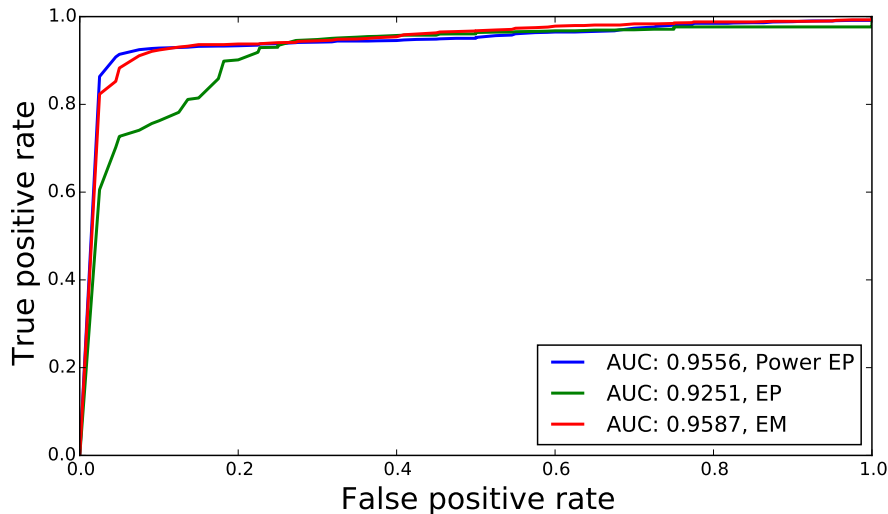
Other training methods for the RTM

- Approximate Expectation Propagation (EP)
- Coordinate-wise optimization of variational evidence lower bound (EM)

Experiments

Synthetic data

ROC curve for feature selection task:



Experiments

Sentiment analysis

Test set classification accuracy:

	Power EP	EM	EP
Dataset 1:	0.9708	0.9659	0.9683
Dataset 2:	0.7523	0.7294	0.7398